# Project 2: Reinforcement Learning

**Shreya S Ramanujam**                                                    SRAMANUJ@STANFORD.EDU

*AA228/CS238, Stanford University*

## 1. Algorithm Descriptions

### 1.1 Small Data Set and Large Data Set

For both the small as well as the large datasets, I used standard Q-learning, where the Q-table is updated directly from observed transitions given in the dataset $(s, a, r, sp)$ and the final policy is greedy in Q-value with respect to all the actions at a particular state. However, for the larger datasets, many states were never visited in the given offline dataset. This was making direct Q-learning unstable, since using the Q-values of unvisited states from the table can cause unwanted error propagations in the Q value estimations.

   To handle this, I modified the algorithm to interpolate Q-values for unseen states using their nearest visited counterparts. Specifically, I maintained a `visited_state` list to keep track of visited states, implying their Q-values are reliable. When updating or querying the Q-table for an unvisited state $s$, the algorithm locates the closest visited state based on state index (between 1 and 100 for small and 1 to 302020 for large) by doing a binary search on this `visited_state` list and uses its Q-values as a proxy in the Bellman update equation. This allows smoother propagation of information across sparsely sampled regions of the state space and a better overall policy. I used this interpolation both in the Bellman update equation (for both $Q(sp)$ and $Q(s)$ where needed) and also for the final greedy policy itself. The pseudocode is given in Algorithm 1.1.

   **Performance Characteristics**: I ran the small dataset for **10** epochs, taking **4.44 seconds**. With small, even 10 epochs were sufficient to give a good policy. I ran the large dataset for **50** epochs, taking **46.53 seconds**. For large, I needed more epochs to arrive at a reasonable policy and estimate for Q, hence the larger number of epochs.

### 1.2 Medium Data Set

For the Medium Data Set, I tried the same algorithm detailed in 1.1, but I was never able to pass the autograder test. I also tried to include another clause in the Q-learning algorithm: if $sp$ is a terminal state ($\geq 0.45$ in the original OpenAI Gymnasium documentation, which becomes State $\geq 457$ after discretization to 500 states), then the target is just $r$, there is no additional $\gamma \max_a Q_{target}(a)$ term. However, this was not able to match the baseline either.

   Finally, I settled on a degree 5 polynomial ridge regression to fit Q values: Poly(X, y). The predictor (X) values for the regression were the normalized [position, velocity] tuples and the targets were the Bellman update targets: $y = r + \gamma \max_a Q(sp, a)$. After fitting the polynomial, we use that to fill up our Q-table and greedily select actions for every state. The pseudocode is given in Algorithm 1.2.

   **Performance Characteristics**: I ran this algorithm for **25** epochs, which took **27.98 seconds**. Any lesser number of epochs did not give a good policy (because the Q function was not well-estimated), which is why I chose to run for 25 epochs.

---

**Algorithm 1** Q-Learning with Nearest-Visited-State Interpolation

---

1: Initialize $Q(s, a) \leftarrow 0$ for all states $s$ and actions $a$
2: Initialize visited states list as empty
3: **for** each epoch **do**
4:     **for** each transition $(s, a, r, sp)$ in shuffled data **do**
5:         **if** $s$ not visited **then**
6:             Mark $s$ as visited
7:             $Q_{update} \leftarrow Q(\text{nearest visited state}, .)$
8:         **else**
9:             $Q_{update} \leftarrow Q(s, .)$
10:         **end if**
11:         **if** $sp$ not visited **then**
12:             $Q(sp, \cdot) \leftarrow Q(\text{nearest visited state}, .)$
13:         **else**
14:             $Q_{target} \leftarrow Q(sp, .)$
15:         **end if**
16:         Compute target: $y \leftarrow r + \gamma \max_a Q_{target}(a)$
17:         Update: $Q(s, a) \leftarrow Q_{update}(a) + \alpha(y - Q_{update}(a))$
18:     **end for**
19: **end for**
20: **for** each state $s$ **do**
21:     **if** $s$ visited **then**
22:         $\pi(s) \leftarrow \arg\max_a Q(s, a)$
23:     **else**
24:         $\pi(s) \leftarrow \arg\max_a Q(\text{nearest visited state}, a)$
25:     **end if**
26: **end for**
27: **return** policy $\pi$

---

**Algorithm 2** Polynomial Regression to find Q Values

---

1: Initialize $Q(s, a) \leftarrow 0$ for all states $s$ and actions $a$
2: **for** each iteration **do**
3:     **for** each transition $(s, a, r, s')$ in dataset **do**
4:         Compute target: $y \leftarrow r + \gamma \max_{a'} Q(s', a')$
5:     **end for**
6:     Fit polynomial regression $Q(s, a) \approx y$ using $[s/n_{\text{states}}, a/n_{\text{actions}}]$ as features
7:     Predict $Q(s, a)$ for all state-action pairs and update $Q$
8: **end for**
9: Return policy: $\pi(s) \leftarrow \arg\max_a Q(s, a)$

---

## 2. Code

```python
import numpy as np
import pandas as pd
import torch.nn as nn
import argparse
import bisect
from sklearn.preprocessing import PolynomialFeatures
from sklearn.pipeline import make_pipeline
from sklearn.linear_model import Ridge
from time import time

#### Helper functions, but I did not end up using them ####
def state_to_indices(position, velocity):
    pos_idx = int((position + 1.2) / (0.6 + 1.2) * 499)  # scale to [0,
    ↪    499]
    vel_idx = int((velocity + 0.07) / (0.07 + 0.07) * 99)  # scale to [0,
    ↪    99]
    return pos_idx, vel_idx

def indices_to_state(pos_idx, vel_idx):
    position = pos_idx / 499 * (0.6 + 1.2) - 1.2
    velocity = vel_idx / 99 * (0.07 + 0.07) - 0.07
    return position, velocity

def state_to_discretized_vals(state):
    pos_idx = (state - 1) % 500
    vel_idx = (state - 1) // 500
    return (pos_idx, vel_idx)

def nearest_visited_state_mountaincar(visited_list, s):
    pos_s , vel_s = state_to_discretized_vals(s + 1)
    best_dist = float('inf')
    nearest = None
    for v in visited_list:
        pos_v , vel_v = state_to_discretized_vals(v + 1)
        dist = abs(pos_v - pos_s) + abs(vel_v - vel_s)
        if dist < best_dist:
            best_dist = dist
            nearest = v
    return nearest

#### Unused Helper functions end ####

# function to find nearest visited neighbour for small and large
```

```python
42  def nearest_visited_state(visited_list, state):
43      # binary search for nearest visited state
44      idx = bisect.bisect_left(visited_list, state)
45      if idx == 0:
46          nearest = visited_list[0]
47      elif idx == len(visited_list):
48          nearest = visited_list[-1]
49      else:
50          left = visited_list[idx - 1]
51          right = visited_list[idx]
52          nearest = left if abs(left - state) <= abs(right - state) else
            ↪  right
53      return nearest
54
55  # code for medium dataset algorithm
56  def fitted_q_iteration_poly(data, n_states, n_actions, gamma=0.99,
    ↪  n_iters=50, degree=5):
57      model = make_pipeline(PolynomialFeatures(degree), Ridge(alpha=1.0))
58      q_values = np.zeros((n_states, n_actions))
59
60      for it in range(n_iters):
61          X, y = [], []
62          for _, row in data.iterrows():
63              s, a, r, sp = row['s'] - 1, row['a'] - 1, row['r'], row['sp'] -
                ↪  1
64              target = r + gamma * np.max(q_values[sp])
65              X.append([s / n_states, a / n_actions]) # to keep range same
66              y.append(target)
67
68          X, y = np.array(X), np.array(y)
69          model.fit(X, y)
70
71          sa_pairs = np.array([[s/n_states, a/n_actions] for s in
            ↪  range(n_states) for a in range(n_actions)])
72          q_pred = model.predict(sa_pairs)
73          q_values = q_pred.reshape(n_states, n_actions)
74          print(f"Iteration {it+1}/{n_iters} done.")
75
76      return np.argmax(q_values, axis=1) + 1
77
78  # code for small and large dataset algorithm
79  def q_learning(task_name, data, n_states, n_actions, alpha=0.1, gamma=0.95,
    ↪  epochs=10):
80      # Initialize Q-table
81      q_table = np.zeros((n_states, n_actions))  # n_states and n_actions
82      # maintain a set of visited states
```

```
83      visited_states = np.zeros(n_states, dtype=bool)
84      visited_list = []   # keep sorted list of visited states
85      for epoch in range(epochs):
86          print(f"Epoch {epoch+1}/{epochs}")
87          for _, row in data.sample(frac=1, random_state=epoch).iterrows():
88              s = row['s'] - 1 # since states are 1-indexed in the data
89              a = row['a'] - 1 # since actions are 1-indexed in the data
90              r = row['r']
91              s_next = row['sp'] - 1 # since states are 1-indexed in the data
92              # for q_table, find nearest state that we have visited
93              # mark visited
94              if not visited_states[s]:
95                  bisect.insert(visited_list, s)  # insert while keeping
                    ↪  sorted
96                  visited_states[s] = True
97                  nearest = nearest_visited_state(visited_list, s)
98                  q_update = q_table[nearest]
99              else:
100                 q_update = q_table[s]
101
102             if not visited_states[s_next]:
103                 bisect.insert(visited_list, s_next)  # insert while keeping
                    ↪  sorted
104                 nearest = nearest_visited_state(visited_list, s_next)
105                 q_target = q_table[nearest]
106             else:
107                 q_target = q_table[s_next]
108
109             # if (task_name == 'medium' and state_to_discretized_vals(s +
                ↪  1)[0] >= 457): # terminal state
110             #     target = r
111             # else:
112             target = r + gamma * np.max(q_target)
113
114             q_table[s, a] = q_update[a] + alpha * (target - q_update[a])
115
116     policy = np.zeros(n_states, dtype=int)
117
118     # Derive policy using nearest visited state via binary search
119     for i in range(n_states):
120         if visited_states[i]:
121             policy[i] = np.argmax(q_table[i])
122         else:
123             nearest = nearest_visited_state(visited_list, i)
124             policy[i] = np.argmax(q_table[nearest])
125
```

```python
126        return policy + 1  # convert back to 1-indexed
127
128
129  def main():
130      # get arguments from command line with defaults
131      parser = argparse.ArgumentParser()
132      parser.add_argument('--task_name', type=str, default='small',
         ↪  help='Task name: small, medium, large')
133      parser.add_argument('--alpha', type=float, default=0.1, help='Learning
         ↪  rate')
134      parser.add_argument('--gamma', type=float, default=None, help='Discount
         ↪  factor')
135      parser.add_argument('--num_epochs', type=int, default=None,
         ↪  help='Number of training epochs')
136      parser.add_argument('--degree', type=int, default=None, help='Degree of
         ↪  polynomial features for medium task')
137      args = parser.parse_args()
138
139      if(args.task_name == 'small'):
140          n_states = 100
141          n_actions = 4
142          input_data_path = './data/small.csv'
143          output_path = 'small.policy'
144          args.gamma = 0.95 if args.gamma is None else args.gamma
145          args.num_epochs = 10 if args.num_epochs is None else
             ↪  args.num_epochs
146
147      elif(args.task_name == 'medium'):
148          n_states = 50000
149          n_actions = 7
150          input_data_path = './data/medium.csv'
151          output_path = 'medium.policy'
152          args.gamma = 1.0 if args.gamma is None else args.gamma
153          args.num_epochs = 25 if args.num_epochs is None else
             ↪  args.num_epochs
154          args.degree = 5 if args.degree is None else args.degree
155
156      elif(args.task_name == 'large'):
157          n_states = 302020
158          n_actions = 9
159          input_data_path = './data/large.csv'
160          output_path = 'large.policy'
161          args.gamma = 0.95 if args.gamma is None else args.gamma
162          args.num_epochs = 50 if args.num_epochs is None else
             ↪  args.num_epochs
163
```

```python
164        # read data
165        data = pd.read_csv(input_data_path)
166        # convert all data to ints
167        data = data.astype(int)
168        if args.task_name == 'medium':
169            print("Using specialized Q-learning for medium task")
170            start_time = time()
171            policy = fitted_q_iteration_poly(data, n_states, n_actions,
               ↪  gamma=args.gamma, n_iters=args.num_epochs, degree=args.degree)
172            print(f"Medium task completed in {time() - start_time:.2f}
               ↪  seconds")
173        else:
174            print("Using general Q-learning")
175            start_time = time()
176            policy = q_learning(args.task_name, data, n_states, n_actions,
               ↪  alpha=args.alpha, gamma=args.gamma, epochs=args.num_epochs)
177            print(f"General Q-learning on {args.task_name} completed in {time()
               ↪  - start_time:.2f} seconds")
178        # output to a file
179        with open(output_path, 'w') as f:
180            for _, action in enumerate(policy):
181                f.write(f"{action}\n")
182
183  if __name__ == "__main__":
184      main()
185
```